

Date	
Page	

- ② Now find the $TF(t, d)$ matrix

4m - (1) sun rising sun. west
(2) sun rising east east

	sum	nbring	east	west.
d ₁	$2/4 = 0.5$	$1/4 = 0.25$	$0/4 = 0$	$1/4 = 0.25$
d ₂	$1/4 = 0.25$	$1/4 = 0.25$	$2/4 = 0.5$	$0/4 = 0$
PDF	$\log\left(\frac{2}{2}\right) = 0$	$\log\left(\frac{2}{2}\right) = 0$	$\log\left(\frac{2}{7}\right) = -0.3$	$\log\left(\frac{2}{7}\right) = -0.3$

→ Now I have to calculate the importance of each keyword in that ^{each} corresponding document

③ Calculate the TF-IDF matrix

→ TF-IDF (t, d, D)

keyword document corpus

→ % TF-IDF = 0 for a word in a document, it means that keyword has no importance in defining that document

	sun	rising	east	west
d1	$0.5 \times 0 = 0$	$0.25 \times 0 = 0$	0	0.075
d2	$0.25 \times 0 = 0$	$0.25 \times 0 = 0$	0.15	0
IDF	0	0	0.3	0.3

Q) How to calculate the similarity between query string and a document??

Ans) Take the dot product of their vectors.

Ex: - Query string: sun (shining) sun west

④ Calculate the TFIDF matrix for query string

	sun	rising	east	west
query	$\frac{2}{3}$	$\frac{0}{3}$	$\frac{0}{3}$	$\frac{1}{3}$

Here as our document has no keyword "shining",

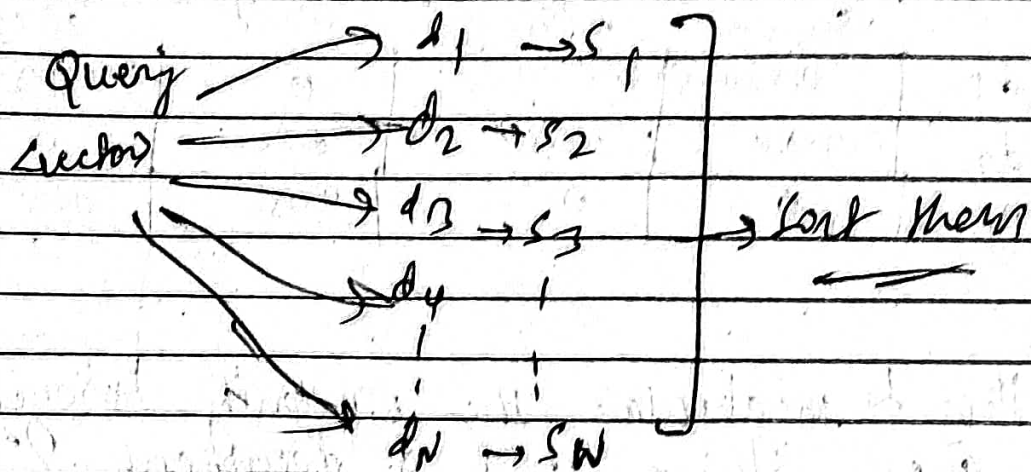
so, we ignore it and consider it as a stopword.

TFIDF

	sun	rising	east	west
query	$\frac{2}{3} \times 0 = 0$	$\frac{0}{3} \times 0 = 0$	$0 \times 0.3 = 0$	$0.3 \times 1 = 0.3$

→ Now we have vectors for each document and a vector for the string.

→ We will calculate the similarity values for that query string with each of the documents.



→ The document having higher similarity values are to be placed on top of the search results.

(5) Find the cosine similarities as explained above.

$$S_i = \frac{\vec{v}_q \cdot \vec{d}_i}{|\vec{v}_q| |\vec{d}_i|}$$

Similarity of i^{th} document with query string

as $\cos \theta = \frac{\vec{v}_1 \cdot \vec{v}_2}{|\vec{v}_1| |\vec{v}_2|}$

explains the inclination / projection.

Ques:-

$$d_1 = \langle 1, 2, 5, 2 \rangle$$

$$q = \langle 3, 0, 4, 1 \rangle$$

$$S_1 = 1 \times 3 + 2 \times 0 + 5 \times 4 + 2 \times 1$$

$$\sqrt{1^2 + 2^2 + 5^2 + 2^2} \sqrt{3^2 + 0^2 + 4^2 + 1^2}$$

→ Store each of the vector (TF-IDF) in their own ~~document~~ separate (i.e.) documents

